

Mixing Time Estimation in Ergodic Markov Chains from a Single Trajectory with Contraction Methods

Geoffrey Wolfer

GEOFFREY@POST.BGU.AC.IL

*Department of Computer Science
Ben-Gurion University of the Negev
Israel*

Editors: Aryeh Kontorovich and Gergely Neu

Abstract

The mixing time t_{mix} of an ergodic Markov chain measures the rate of convergence towards its stationary distribution π . We consider the problem of estimating t_{mix} from one single trajectory of m observations (X_1, \dots, X_m) , in the case where the transition kernel M is unknown, a research program started by [Hsu et al. \(2015\)](#). The community has so far focused primarily on leveraging spectral methods to estimate the *relaxation time* t_{rel} of a *reversible* Markov chain as a proxy for t_{mix} . Although these techniques have recently been extended to tackle non-reversible chains, this general setting remains much less understood. Our new approach based on contraction methods is the first that aims at directly estimating t_{mix} up to multiplicative small universal constants instead of t_{rel} . It does so by introducing a generalized version of Dobrushin’s contraction coefficient κ_{gen} , which is shown to control the mixing time regardless of reversibility. We subsequently design fully data-dependent high confidence intervals around κ_{gen} that generally yield better convergence guarantees and are more practical than state-of-the-art.

Keywords: Ergodic Markov chain, mixing time, Dobrushin contraction coefficient

1. Introduction

The topic of this work is the construction of a non-trivial high confidence interval around the mixing time of a finite state ergodic Markov chain, when one is only allowed to observe a single long trajectory of states X_1, X_2, \dots, X_m , i.e. does not have access to a restart mechanism. The problem is motivated by PAC-type learning problems that assume data sampled from a Markovian process, where generalization guarantees oftentimes involve the *a priori* unknown mixing properties of the chain. Other applications are in MCMC diagnostics for non-reversible Markov chains, that may enjoy better mixing properties or asymptotic variance than their reversible counterparts, or in the field of reinforcement learning, where bounds on the mixing time are routinely assumed. We invite the reader to the *related work* sections of [Hsu et al. \(2019\)](#); [Wolfer and Kontorovich \(2019\)](#) for a complete set of references to the aforementioned problems and additional motivation.

Main contributions.

- In Section 2, in lieu of the (pseudo-)relaxation time t_{rel} , we propose a new proxy for the mixing time based on a *contraction coefficient* κ_{gen} that generalizes Dobrushin’s, and in particular, does not require reversibility. We show in Theorem 1 that this quantity controls the mixing time up to multiplicative universal constants – which are small and given at (10) –

such that contrary to the relaxation time, it is not subject to the gap mentioned at (5). Namely,

$$\frac{1}{1 - \kappa_{\text{gen}}} = \Theta(t_{\text{mix}}).$$

- In Section 3.1, we design fully empirical confidence intervals around κ_{gen} that in the general (non-reversible) setting are thinner, and considerably more practical than their spectral state-of-the-art counterparts: For a chain on d states and a chosen parameter $S \in \mathbb{N}$, our estimator $\widehat{\kappa}_{\text{gen}[S]}$ is such that

$$|\widehat{\kappa}_{\text{gen}[S]} - \kappa_{\text{gen}}| \leq \tilde{\mathcal{O}} \left(\frac{1}{S} + \max_{\ell \in [S]} \left\{ \frac{1}{\ell} \sqrt{\frac{d}{N_{\text{min}}^{(\ell)}}} \right\} \right),$$

where $N_{\text{min}}^{(\ell)}$ is the *least number of visits for the ℓ -skipped chain*, a fully observable quantity defined at (11). Additionally, the analysis leading to the confidence intervals is of an arguably much simpler nature than that of [Wolfer and Kontorovich \(2019\)](#).

- In Section 3.2, for a d state Markov chain with minimum stationary probability π_{\star} (definition at (2)), we further deduce point estimators for estimating κ_{gen} down to absolute error ε (Theorem 5), with sample complexity $m_{+} = \tilde{\mathcal{O}} \left(\frac{1}{\pi_{\star}} \max \left\{ t_{\text{mix}}, \frac{d}{\varepsilon^2} \right\} \right)$ and relative error ε (Theorem 6), for a trajectory length of $m_{\times} = \tilde{\mathcal{O}} \left(\frac{dt_{\text{mix}}^2}{\pi_{\star} \varepsilon^2} \right)$, offering better guarantees than the one of [Wolfer and Kontorovich \(2019\)](#) for the non-trivial classes of slow mixing chains ($t_{\text{mix}} > d$), and chains whose stationary distribution is not close to being uniform (see Remark 8).

Notation and setting. The set \mathbb{N} will refer to $\{1, 2, 3, \dots\}$ and for $n \in \mathbb{N}$, we write $[n] = \{1, 2, 3, \dots, n\}$. Let Ω a set such that $|\Omega| = d < \infty$, and define Δ_{Ω} the simplex of all distributions – seen as row vectors – over Ω . For $(\boldsymbol{\mu}, \boldsymbol{\nu}) \in \Delta_{\Omega}^2$, we define the *total variation distance* in terms of the ℓ_1 norm:

$$\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_{\text{TV}} \doteq \frac{1}{2} \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_1. \quad (1)$$

We consider *time-homogeneous Markov chains*

$$X_1, X_2, \dots, X_t, \dots \sim (\boldsymbol{\mu}, \mathbf{M})$$

with *initial distribution* $\boldsymbol{\mu} \in \Delta_{\Omega}$, and *row-stochastic transition matrix* $\mathbf{M}: \Omega \times \Omega \rightarrow [0, 1]$. We say that a Markov chain $(\boldsymbol{\mu}, \mathbf{M})$ is *ergodic* when \mathbf{M} is a *primitive matrix*, i.e. $\exists p \in \mathbb{N}, \mathbf{M}^p > 0$ entry-wise. In this case, \mathbf{M} has a unique *stationary distribution* $\boldsymbol{\pi}$ such that $\boldsymbol{\pi} \mathbf{M} = \boldsymbol{\pi}$, the chain is known to converge to $\boldsymbol{\pi}$, and the *minimum stationary probability*

$$\pi_{\star} \doteq \min_{i \in \Omega} \boldsymbol{\pi}(i) \quad (2)$$

is such that $\pi_{\star} > 0$. We measure distance to stationarity in total variation,

$$h(t) \doteq \sup_{\boldsymbol{\mu} \in \Delta_{\Omega}} \|\boldsymbol{\mu} \mathbf{M}^t - \boldsymbol{\pi}\|_{\text{TV}}. \quad (3)$$

For $\xi \in (0, 1/2)$, the *mixing time* of M is defined by

$$t_{\text{mix}}(\xi) \doteq \arg \min_{t \in \mathbb{N}} \{h(t) < \xi\}, \quad (4)$$

and by convention $t_{\text{mix}} \doteq t_{\text{mix}}(1/4)$. The reader is referred to [Levin et al. \(2009, Chapter 4\)](#) for a more detailed introduction to Markov chain mixing. We will use the standard \mathcal{O} and Θ notations, and $\tilde{\mathcal{O}}, \tilde{\Theta}$ when logarithmic dependencies in any natural parameter are omitted, and for $x \in \mathbb{R}_+$, we will use $\tilde{\ln} x$ as a shorthand for $\ln x \ln \ln x$. The definition of elements specific to contraction methods is deferred to [Section 2](#) for a clearer exposition.

Related work. Research has so far mostly focused on leveraging spectral methods for estimating the *relaxation time* t_{rel} of a *reversible* ([Levin et al., 2009, Section 1.6](#)) chain as an approximation of t_{mix} ([Hsu et al., 2015](#); [Levin and Peres, 2016](#); [Hsu et al., 2019](#); [Combes and Touati, 2019](#); [Qin et al., 2019](#)). Indeed, in this setting, t_{rel} , defined as the inverse of the *absolute spectral gap* γ_* , is known to be related – see [\(5\)](#) – to the mixing time up to a logarithmic correction ([Levin et al., 2009, Theorem 12.4](#)). The problem, therefore, reduces to estimating the second largest eigenvalue in magnitude. Moreover, as reversibility and self-adjointness of the Markov operator are equivalent notions, dimension-free perturbation eigenvalue bounds – namely Weyl’s inequality – are available to efficiently estimate its spectrum, which is the subject of [Hsu et al. \(2015\)](#); [Levin and Peres \(2016\)](#); [Hsu et al. \(2019\)](#). In their work [Combes and Touati \(2019\)](#), offer a different perspective on the problem by putting the emphasis on computational complexity, invoking *power methods* and *upper confidence interval* techniques to design a more space-efficient estimator. Finally, [Qin et al. \(2019\)](#) explore the case of general state spaces, for kernels that are *trace-class* operators (compact with summable eigenvalues). Although broad collections of chains are known to be reversible such as random walks on graphs or birth and death processes, this assumption is a strong restriction on the class of chains that can be treated, and our approach compares favorably with this body of work in as much as it removes this requirement.

One exception to the above list is [Wolfer and Kontorovich \(2019\)](#) that extends the estimation results to the *non-reversible* setting, by estimating the pseudo-relaxation time ([Kamath and Verdú, 2016, \(16\)](#)). More specifically, even in the absence of reversibility, it was shown in [Paulin \(2015, Proposition 3.4\)](#) that a related quantity, the inverse of the *pseudo-spectral gap*

$$\gamma_{\text{ps}} \doteq \max_{k \in \mathbb{N}} \left\{ \gamma((M^\dagger)^k M^k) / k \right\},$$

where M^\dagger is the time-reversal of M , still traps the mixing time up to a logarithmic correction. [Wolfer and Kontorovich \(2019\)](#) carry out the analysis of estimating this quantity, with a scheme that consists in observing multi-step chains forward and backward in time. They show that it is enough to explore a finite set of skipping rates, and recover a consistent estimator that still enjoys spectral stability, and converges to arbitrary precision with a trajectory length polynomial in the natural parameters $d, \pi_*, \varepsilon, t_{\text{mix}}$.

An inherent drawback of *all* previous approaches, however, is the existence of a known gap between the (pseudo-)relaxation time and t_{mix} that depends on π_* or d ([Levin et al., 2009, Theorem 12.4](#)), ([Paulin, 2015, Proposition 3.4](#)), ([Jerison, 2013, Theorem 1.2](#)). We can summarize these results as

$$c_1 \cdot (t_{\text{rel}} - 1) \leq t_{\text{mix}} \leq c_2 \cdot \min \left\{ d, \ln \frac{1}{\pi_*} \right\} t_{\text{rel}}, \text{ where } (c_1, c_2) \in \mathbb{R}_+^2. \quad (5)$$

Moreover, it is known that this gap cannot generally be closed (Jerison, 2013), so that estimation of t_{rel} down to arbitrary error still will not yield an accurate estimate for t_{mix} as $d \rightarrow \infty$. This limitation is the motivation behind our search for a new, tighter proxy. Although the task of estimating t_{mix} directly is currently believed to be more challenging than t_{rel} as raised in the concluding remarks of Combes and Touati (2019, Conclusion), no rigorous or quantitative comparison is known in terms of statistical complexity. This work therefore also initiates the investigations towards answering this question. Comparisons of convergence rates with the state-of-the-art point empirical confidence intervals and estimators are respectively carried out at Remark 4 and Remark 8.

Finally, as raised in Combes and Touati (2019), there exists an interesting trade-off between computational complexity and statistical accuracy. This work is more concerned with the latter, such that the designed procedures will be applicable for medium-sized state spaces, with computational complexities of the same order to that of Wolfer and Kontorovich (2019).

2. Generalized contraction coefficient

For a Markov chain M , the *Dobrushin contraction coefficient*, also known as *Dobrushin ergodic coefficient* (Dobrushin, 1956), (Brémaud, 1999, Definition 7.1) is defined by

$$\kappa \doteq \max_{(i,j) \in \Omega^2} \|M(i, \cdot) - M(j, \cdot)\|_{\text{TV}}, \quad (6)$$

where the term *contraction* refers to the property (Brémaud, 1999, Corollary 7.1) that $\forall(\boldsymbol{\mu}, \boldsymbol{\nu}) \in \Delta_{\Omega}^2$,

$$\|(\boldsymbol{\mu} - \boldsymbol{\nu})M\|_{\text{TV}} \leq \kappa \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_{\text{TV}}. \quad (7)$$

Contraction in the sense of Dobrushin is a special case of *coarse Ricci curvature* (Ollivier, 2009), where the metric taken on Ω is the *discrete metric*, and the *Wasserstein distance* between distributions reduces to total variation. In the case where $\kappa < 1$, the Bubley-Dyer path coupling bound (Bubley and Dyer, 1997) gives an upper bound on mixing time

$$t_{\text{mix}}(\xi) \leq \frac{\ln \xi}{\ln(1 - \kappa)}$$

Unfortunately, there exists a large subset of ergodic chains such that $\kappa = 1$, and for which this direct method fails to yield convergence rates. To overcome this limitation, we consider multi-step chains, where for $s \in [m - 1]$,

$$X_1, X_{1+s}, X_{1+2s}, \dots, X_{1+\lfloor(m-1)/s\rfloor s} \sim (\boldsymbol{\mu}, M^s),$$

and define the contraction coefficient of the chain with skipping rate s to be

$$\kappa_s \doteq \max_{(i,j) \in \Omega^2} \|M^s(i, \cdot) - M^s(j, \cdot)\|_{\text{TV}}. \quad (8)$$

We then introduce a *generalized contraction coefficient* κ_{gen} of the ergodic chain M :

$$\kappa_{\text{gen}} \doteq 1 - \max_{s \in \mathbb{N}} \left\{ \frac{1 - \kappa_s}{s} \right\}, \quad (9)$$

and write s_{gen} the smallest integer ¹ such that $\kappa_{\text{gen}} = 1 - \frac{1 - \kappa_{s_{\text{gen}}}}{s_{\text{gen}}}$. This quantity is derived in a similar spirit as [Paulin \(2015\)](#) defined the pseudo-spectral gap of an ergodic Markov chain. Even in the case where $\kappa = 1$, we now formalize in [Theorem 1](#) the fact that $\kappa_{\text{gen}} < 1$ always holds, as $\frac{1}{1 - \kappa_{\text{gen}}}$ traps t_{mix} up to universal constants.

Theorem 1 *Let $\xi \in (0, 1/2)$, and M ergodic with mixing time $t_{\text{mix}}(\xi)$, then*

$$\frac{1 - 2\xi}{1 - \kappa_{\text{gen}}} \leq t_{\text{mix}}(\xi) \leq \frac{1 + \ln 1/\xi}{1 - \kappa_{\text{gen}}},$$

where κ_{gen} is defined at [\(9\)](#), and in particular,

$$\frac{1/2}{1 - \kappa_{\text{gen}}} \leq t_{\text{mix}} \leq \frac{1 + \ln 4}{1 - \kappa_{\text{gen}}}. \quad (10)$$

We point out that although we could not find any reference to the quantity at [\(9\)](#), or to [Theorem 1](#), considering multi-step contractions to study concentration or mixing properties of chains is not a novel idea in itself; see for instance [Dyer et al. \(2001\)](#); [Luczak \(2008\)](#); [Paulin \(2016\)](#).

3. Statistical estimation of the mixing time from a single trajectory

This section is devoted to the analysis of the statistical complexity of estimating the mixing time of an ergodic chain from one single long draw of observations (no restart mechanism). [Section 2](#) introduced κ_{gen} [defined at [\(9\)](#)] as a tighter proxy for t_{mix} [as shown by [Theorem 1](#)] and allows for a reduction of the estimation problem. In [Section 3.1](#) we construct fully empirical high-confidence intervals around κ_{gen} . We further derive point estimators in [Section 3.2](#) and analyze their finite sample convergence properties both in *absolute* ([Theorem 5](#)) and *relative* ([Theorem 6](#)) error.

3.1. Fully empirical confidence intervals

For a confidence parameter δ , and a trajectory X_1, \dots, X_m , our goal is to construct a non-trivial interval $I_{\kappa_{\text{gen}}} = (\kappa_{\text{gen,lb}}, \kappa_{\text{gen,ub}})$ such that

$$\mathbb{P}(\kappa_{\text{gen}} \in I_{\kappa_{\text{gen}}}) \geq 1 - \delta.$$

Our estimator will be a truncated plug-in version of κ_{gen} , where we only explore a prefix $[S]$ of the integers, a similar idea as employed in [Wolfer and Kontorovich \(2019\)](#) for estimating the pseudo-spectral gap ([Paulin, 2015](#), [Section 3.1](#)). For a chain with skipping rate s , we define the following random variables,

$$\begin{aligned} N_i^{(s)} &\doteq \sum_{t=1}^{\lfloor (m-1)/s \rfloor} \mathbf{1}\{X_{1+s(t-1)} = i\}, & N_{\min}^{(s)} &\doteq \min_{i \in \Omega} N_i^{(s)}, \\ N_{ij}^{(s)} &\doteq \sum_{t=1}^{\lfloor (m-1)/s \rfloor} \mathbf{1}\{X_{1+s(t-1)} = i, X_{1+st} = j\}, \end{aligned} \quad (11)$$

1. The existence of s_{gen} is guaranteed by the observation that $s \mapsto \frac{1 - \kappa_s}{s} \in (0, \frac{1}{s})$.

and construct an estimator for the multi-step kernel M^s and its contraction coefficient,

$$\widehat{M}^{(s)} \doteq \sum_{(i,j) \in \Omega^2} \frac{N_{ij}^{(s)}}{N_i^{(s)}} \mathbf{1}\{N_i^{(s)} > 0\} \mathbf{e}_i \otimes \mathbf{e}_j, \quad \widehat{\kappa}_s \doteq \kappa \left(\widehat{M}^{(s)} \right), \quad (12)$$

where \mathbf{e}_i is the i th coordinate basis vector and \otimes denotes the standard tensor product. When $s = 1$, we will omit subscript or superscript and write respectively $N_i, N_{ij}, N_{\min}, \widehat{M}$ and $\widehat{\kappa}$. Finally, the estimator for κ_{gen} parametrized by an integer S is

$$\widehat{\kappa}_{\text{gen}[S]}: \Omega^m \rightarrow (0, 1), \quad \mathbf{X} \mapsto 1 - \max_{s \in [S]} \left\{ \frac{1 - \widehat{\kappa}_s(\mathbf{X})}{s} \right\}.$$

Theorem 2 *Let $\delta \in (0, 1)$, $S \in \mathbb{N}$, and $X_1, \dots, X_m \sim M$, then with probability $1 - \delta$,*

$$|\widehat{\kappa}_{\text{gen}[S]} - \kappa_{\text{gen}}| \leq \frac{1}{S} + \sqrt{d} \max_{s \in [S]} \left\{ \frac{\mathcal{L}_s}{s \sqrt{N_{\min}^{(s)}}} \right\},$$

where $\mathcal{L}_s = \mathcal{O} \left(\ln \left(\frac{dS \ln m/s}{\delta} \right) \right)$, and $N_{\min}^{(s)}$ is defined at (11).

Remark 3 *As we choose to carry out our analysis with unsmoothed estimators, we see that the confidence bounds can be ill-defined for short trajectories. A slight modification of the proofs with a smoothing parameter; for example, analyzing*

$$\widehat{M}^{(\lambda, s)} \doteq \sum_{(i,j) \in \Omega^2} \frac{N_{ij}^{(s)} + \lambda}{N_i^{(s)} + d\lambda} \mathbf{e}_i \otimes \mathbf{e}_j,$$

with $\lambda > 0$ instead, can yield intervals that are well defined almost surely. This would, however, clutter the analysis while offering only incremental improvement.

Remark 4 *Not only is the interval at Theorem 2 far more user-friendly than the one designed around the pseudo-spectral gap in Wolfer and Kontorovich (2019, Theorem 8), it is also much narrower. Denoting by \approx a rough estimate of the rate at which we expect the intervals to decay in width,*

$$|\widehat{\kappa}_{\text{gen}[S]} - \kappa_{\text{gen}}| \approx \frac{1}{S} + \sqrt{\frac{d}{\pi_* m}}, \quad (13)$$

whereas the known intervals around the pseudo-spectral gap γ_{ps} of the estimator $\widehat{\gamma}_{\text{ps}[S]}$ defined by Wolfer and Kontorovich could generally decay as slowly as

$$|\widehat{\gamma}_{\text{ps}[S]} - \gamma_{\text{ps}}| \approx \frac{1}{S} + \frac{1}{\pi_*^{3/2}} \sqrt{\frac{d}{m}} \left(\sqrt{d} + \frac{1}{\gamma_{\text{ps}}} \right).$$

We end this section with a short discussion on the choice of S , that is missing from Wolfer and Kontorovich (2019). Assuming a $1 - \delta/2$ confidence interval $I_{\pi_*} = (\pi_{*,\text{lb}}, \pi_{*,\text{ub}})$ around π_* , for instance employing the estimation procedure of Hsu et al. (2015), then a practical choice of S for balancing the two terms at (13) is

$$S \approx n \vee \sqrt{m(\pi_{*,\text{lb}} \wedge 1/d)/d}, \quad (14)$$

where n is a small arbitrary integer. In other words, it is reasonable to wait for the trajectory length to be of the order of the square of the state space size before starting to explore larger skipping rates.

3.2. Point estimator for κ_{gen}

For chosen precision ε and confidence δ parameters, we construct a point estimator down to absolute error, where the algorithm only needs knowledge of d and ε, δ in order to run.

Theorem 5 (Point estimator for κ_{gen} (absolute error)) *Let $(\varepsilon, \delta) \in (0, 1)^2$, and let*

$$X_1, X_2, \dots, X_m \sim (\boldsymbol{\mu}, \mathbf{M})$$

an unknown ergodic Markov chain with minimum stationary probability π_ , and generalized contraction coefficient κ_{gen} . There exists an estimation procedure $\widehat{\kappa}_{\text{gen}}^+$: $\Omega^m \rightarrow (0, 1)$ such that for*

$$m \geq c \frac{\mathcal{L}}{\pi_*} \max \left\{ \frac{1}{1 - \kappa_{\text{gen}}}, \frac{d}{\varepsilon^2} \right\},$$

$|\widehat{\kappa}_{\text{gen}}^+ - \kappa_{\text{gen}}| < \varepsilon$ holds with probability at least $1 - \delta$, where $\mathcal{L} = \mathcal{O} \left(\widetilde{\ln} \left(\frac{d}{\delta \varepsilon \pi_} \right) \right)$, and c is a universal constant.*

From the proof of Theorem 5, we observe that, perhaps surprisingly, for a contracting chain, i.e. $\kappa = 1 - \alpha$ with $\alpha > 0$, the statistical difficulty of estimating α is of the same order (ignoring logarithmic factors) as that of estimating κ_{gen} , while only providing with – generally sub-optimal– upper bounds on the mixing time. One remaining question is the necessity of the dependency in d . A heuristic argument based on the results of Jiao et al. (2018) would seem to imply that any technique basing itself solely on the definition of a contraction coefficient – boiling down to estimating ℓ_1 distances of $[d]$ supported distributions – would necessarily have a statistical dependency in the support size.

We now show that it is also possible to construct an algorithm that outputs an estimate of $1 - \kappa_{\text{gen}}$ with relative error ε . For this goal, the algorithm needs to explore at least the first $S = \Theta \left(\frac{1}{\varepsilon(1 - \kappa_{\text{gen}})} \right)$, involving the unknown quantity κ_{gen} . The solution is an adaptive argument that is fleshed out in Section 4.4.

Theorem 6 (Point estimator for $1 - \kappa_{\text{gen}}$ (relative error)) *Let $(\varepsilon, \delta) \in (0, 1)^2$, and let*

$$X_1, X_2, \dots, X_m \sim (\boldsymbol{\mu}, \mathbf{M})$$

an unknown ergodic Markov chain with minimum stationary probability π_ , and generalized contraction coefficient κ_{gen} . There exists an estimation procedure $\widehat{\kappa}_{\text{gen}}^\times$: $\Omega^m \rightarrow (0, 1)$ such that for*

$$m \geq c \frac{\mathcal{L}d}{\pi_*(1 - \kappa_{\text{gen}})^2 \varepsilon^2},$$

$\left| \frac{1 - \widehat{\kappa}_{\text{gen}}^\times}{1 - \kappa_{\text{gen}}} - 1 \right| < \varepsilon$ holds with probability at least $1 - \delta$, where $\mathcal{L} = \mathcal{O} \left(\widetilde{\ln} \left(\frac{d}{\delta \varepsilon \pi_(1 - \kappa_{\text{gen}})} \right) \right)$, and c is a universal constant.*

Corollary 7 (to Theorem 6) *Let $\delta \in (0, 1)$, and let $X_1, X_2, \dots, X_m \sim (\boldsymbol{\mu}, \mathbf{M})$ an unknown ergodic Markov chain with minimum stationary probability π_* , and mixing time t_{mix} . There exists an estimation procedure $\hat{t}_{\text{mix}}: \Omega^m \rightarrow \mathbb{N}$ such that for $m \geq c \frac{\mathcal{L}d t_{\text{mix}}^2}{\pi_*}$, $\frac{1}{3}t_{\text{mix}} \leq \hat{t}_{\text{mix}} \leq 3t_{\text{mix}}$, holds with probability at least $1 - \delta$, where $\mathcal{L} = \mathcal{O} \left(\widetilde{\ln} \left(\frac{t_{\text{mix}}d}{\delta \pi_*} \right) \right)$, and c is a universal constant.*

Remark 8 *Although in principle, direct comparison of point estimators with prior research is not possible as all previous work focused on t_{rel} , we treat for now the question as if estimation of κ_{gen} or t_{mix} directly is not harder, and focus on relative error. In the general (non-reversible) setting, the only known finite sample upper bound (Wolfer and Kontorovich, 2019, Theorem 3) for estimating the pseudo-relaxation time is of*

$$m_{\times} = \tilde{O} \left(\frac{t_{\text{rel}}^2}{\pi_{\star} \varepsilon^2} \max \{t_{\text{rel}}, \beta(\boldsymbol{\pi}) \min \{\beta(\boldsymbol{\pi}), d\}\} \right),$$

where $\beta(\boldsymbol{\pi}) \doteq \max_{(i,j) \in \Omega^2} \left\{ \frac{\pi(i)}{\pi(j)} \right\}$ measures how far $\boldsymbol{\pi}$ is from being uniform, and $1 \leq \beta(\boldsymbol{\pi}) \leq 1/\pi_{\star}$. From Theorem 6, it is possible to estimate $1 - \kappa_{\text{gen}}$ with $m_{\times} = \tilde{O} \left(\frac{dt_{\text{mix}}^2}{\pi_{\star} \varepsilon^2} \right)$, so that for the two classes of slow mixing chains ($t_{\text{mix}} > d$) and chains with a stationary distribution $\boldsymbol{\pi}$ such that $\beta(\boldsymbol{\pi}) > \sqrt{d}$, our result dominates, showing that the two methods offer complementary convergence rates.

4. Proofs

4.1. Proof of Theorem 1.

The proof is standard. See for example Paulin (2015, Section 5.2) for a similar technique. We will first bound the distance to stationarity $h(t)$ for a given $t > s_{\text{gen}}$,

$$\begin{aligned} h(t) &\doteq \sup_{\boldsymbol{\mu} \in \Delta_{\Omega}} \left\| \boldsymbol{\mu} \mathbf{M}^t - \boldsymbol{\pi} \right\|_{\text{TV}} \\ &\stackrel{(i)}{=} \sup_{\boldsymbol{\mu} \in \Delta_{\Omega}} \left\| (\boldsymbol{\mu} \mathbf{M}^{t-s_{\text{gen}}} - \boldsymbol{\pi}) \mathbf{M}^{s_{\text{gen}}} \right\|_{\text{TV}} \\ &\stackrel{(ii)}{\leq} \kappa(\mathbf{M}^{s_{\text{gen}}}) \sup_{\boldsymbol{\mu} \in \Delta_{\Omega}} \left\| \boldsymbol{\mu} \mathbf{M}^{t-s_{\text{gen}}} - \boldsymbol{\pi} \right\|_{\text{TV}} \\ &\stackrel{(iii)}{\leq} \kappa(\mathbf{M}^{s_{\text{gen}}})^{\lfloor t/s_{\text{gen}} \rfloor} \sup_{\boldsymbol{\mu} \in \Delta_{\Omega}} \left\| \boldsymbol{\mu} \mathbf{M}^{t - \lfloor t/s_{\text{gen}} \rfloor s_{\text{gen}}} - \boldsymbol{\pi} \right\|_{\text{TV}} \\ &\stackrel{(iv)}{\leq} \kappa(\mathbf{M}^{s_{\text{gen}}})^{(t-s_{\text{gen}})/s_{\text{gen}}} \end{aligned}$$

where (i) is by definition of $\boldsymbol{\pi}$, (ii) is the contraction property at (7), (iii) is by an inductive argument, and (iv) is by property of the total variation distance. Since $1 - \kappa_{\text{gen}} = \frac{1 - \kappa(\mathbf{M}^{s_{\text{gen}}})}{s_{\text{gen}}} \leq \frac{1}{s_{\text{gen}}}$, and from properties of the exponential function, $h(t) \leq e \cdot e^{-t(1-\kappa_{\text{gen}})}$, so that for $t > \frac{\ln e/\xi}{1-\kappa_{\text{gen}}}$, $h(t) \leq \xi$, hence the upper bound.

For the lower bound, notice that $\forall (i, j) \in \Omega^2$, by sub-additivity of the ℓ_1 norm and by definition of successively $h(t)$ and $t_{\text{mix}}(\xi)$,

$$\left\| \mathbf{M}^{t_{\text{mix}}(\xi)}(i, \cdot) - \mathbf{M}^{t_{\text{mix}}(\xi)}(j, \cdot) \right\|_{\text{TV}} \leq 2h(t_{\text{mix}}(\xi)) \leq 2\xi,$$

such that by definition of κ_{gen} and κ ,

$$1 - \kappa_{\text{gen}} \geq \frac{1 - \kappa(\mathbf{M}^{t_{\text{mix}}(\xi)})}{t_{\text{mix}}(\xi)} \geq \frac{1 - 2\xi}{t_{\text{mix}}(\xi)}.$$

The lower bound is also a consequence of (Paulin, 2016, Proposition 3.3, (3.2)). □

4.2. Proof of Theorem 2

We first report Wolfer and Kontorovich (2019, Lemma D.4), which we will use in our argument.

Lemma 9 (Wolfer and Kontorovich (2019)) *Let $X_1, \dots, X_m \sim (M, \mu)$ a d -state ergodic Markov chain. Then, with probability at least $1 - \delta$,*

$$\|\widehat{M} - M\|_\infty \leq 4\mathcal{L}\sqrt{\frac{d}{N_{\min}}},$$

where

$$\mathcal{L} \doteq \arg \min_{t \geq 1} \left\{ (1 + \lceil \ln(2m/t) \rceil_+) (d+1)e^{-t} \leq \delta/d \right\} = \mathcal{O} \left(\ln \left(\frac{d \ln m}{\delta} \right) \right),$$

\widehat{M} is the empirical transition matrix of counts defined at (12), and N_{\min} is defined at (11).

In other words, Lemma 9 shows that it is possible to control with high-probability the error in estimating the Markov kernel w.r.t the ℓ_∞ operator norm in terms of the least number of visits. Writing for convenience $\kappa_{\text{gen}[S]} \doteq \max_{s \in [S]} \left\{ \frac{1 - \kappa_s}{s} \right\}$, and for $r \in S$,

$$\mathcal{L}_r \doteq \arg \min_{t \geq 1} \left\{ (1 + \lceil \ln(2m/(tr)) \rceil_+) (d+1)e^{-t} \leq \frac{\delta}{dS} \right\} = \mathcal{O} \left(\ln \left(\frac{dS \ln m/r}{\delta} \right) \right).$$

Then successively,

$$\begin{aligned} & \mathbb{P} \left(\left| \widehat{\kappa}_{\text{gen}[S]} - \kappa_{\text{gen}} \right| > \frac{1}{S} + \max_{r \in [S]} \left\{ \frac{4}{r} \mathcal{L}_r \sqrt{\frac{d}{N_{\min}^{(r)}}} \right\} \right) \\ & \leq \mathbb{P} \left(\max_{s \in [S]} \left| \frac{\widehat{\kappa}_s - \kappa_s}{s} \right| > \max_{r \in [S]} \left\{ \frac{4}{r} \mathcal{L}_r \sqrt{\frac{d}{N_{\min}^{(r)}}} \right\} \right) \\ & \stackrel{(ii)}{\leq} \sum_{s=1}^S \mathbb{P} \left(\left| \widehat{\kappa}_s - \kappa_s \right| > s \max_{r \in [S]} \left\{ \frac{4}{r} \mathcal{L}_r \sqrt{\frac{d}{N_{\min}^{(r)}}} \right\} \right) \\ & \stackrel{(iii)}{\leq} \sum_{s=1}^S \mathbb{P} \left(\left\| \widehat{M}^{(s)} - M^s \right\|_\infty > 4\mathcal{L}_s \sqrt{\frac{d}{N_{\min}^{(s)}}} \right) \\ & \stackrel{(iv)}{\leq} \sum_{s=1}^S \frac{\delta}{S} = \delta, \end{aligned}$$

where (i) follows from the fact that $|\kappa_{\text{gen}} - \kappa_{\text{gen}[S]}| \leq 1/S$ and that for $\nu, \theta \in \mathbb{R}^S$ it is the case from sub-additivity of the uniform norm that $\|\nu\|_\infty - \|\theta\|_\infty \leq \|\nu - \theta\|_\infty$; (ii) is an application of the union bound, (iii) stems from the fact that the ℓ_∞ operator norm dominates the distance between Dobrushin contraction coefficients (Fact 5.1), and (iv) is Lemma 9. □

4.3. Proof of Theorem 5

To reach arbitrary precision down to additive error, we explore the first $S = \lceil 2/\varepsilon \rceil$ possible skipping rates, i.e. consider the estimator $\widehat{\kappa}_{\text{gen}}^{\lceil 2/\varepsilon \rceil}$. Then, following the same first steps (i), (ii), (iii) as in the proof of Theorem 2 together with $S > \frac{2}{\varepsilon}$,

$$\mathbb{P}_\pi \left(\left| \widehat{\kappa}_{\text{gen}}^{\lceil 2/\varepsilon \rceil} - \kappa_{\text{gen}} \right| > \varepsilon \right) \leq \sum_{s=1}^{\lceil 2/\varepsilon \rceil} \mathbb{P}_\pi \left(\left\| \widehat{M}^{(s)} - M^s \right\|_\infty > s \frac{\varepsilon}{2} \right).$$

For each term,

$$\begin{aligned} \mathbb{P}_\pi \left(\left\| \widehat{M}^{(s)} - M^s \right\|_\infty > s \frac{\varepsilon}{2} \right) &\stackrel{(i)}{\leq} \mathbb{P}_\pi \left(\left\| \widehat{M}^{(s)} - M^s \right\|_\infty > 4\mathcal{L}_s \sqrt{\frac{d}{N_{\min}^{(s)}}} \right) \\ &\quad + \mathbb{P}_\pi \left(N_{\min}^{(s)} < \frac{64\mathcal{L}_s^2 d}{s^2 \varepsilon^2} \right) \\ &\stackrel{(ii)}{\leq} \frac{\delta}{2\lceil 2/\varepsilon \rceil} + \mathbb{P}_\pi \left(N_{\min}^{(s)} < \frac{1}{2} \lceil (m-1)/s \rceil \pi_\star \right) \end{aligned}$$

where (i) stems from the fact that for functions of the sample ϕ and ψ , the chaining argument

$$\mathbb{P}_\pi (\phi(\mathbf{X}) > \varepsilon) \leq \mathbb{P}_\pi (\phi(\mathbf{X}) > \psi(\mathbf{X})) + \mathbb{P}_\pi (\psi(\mathbf{X}) > \varepsilon)$$

holds, and (ii) follows from the proof of Theorem 2 at confidence $1 - \delta/2$ for the former summand, and by already setting $m \geq c' \frac{d}{\pi_\star s \varepsilon^2} \widetilde{\ln} \left(\frac{dS}{\delta \pi_\star \varepsilon} \right)$, entailing the sufficient $m \geq \frac{128d\mathcal{L}_s^2}{\pi_\star s \varepsilon^2}$ for the latter. The remaining error probability, which corresponds to an unreasonable number of visits to the least visited state is controlled for $m \geq c'' \left(\frac{1}{\varepsilon} + \frac{t_{\text{mix}}}{\pi_\star} \ln \frac{d}{\varepsilon \delta} \right)$ as a result of Lemma 10, which in turn is a consequence of Chung et al. (2012, Theorem 3.1). Finally, Paulin (2015, Proposition 3.10) extends the bound to non-stationary chains. \square

4.4. Proof of Theorem 6

Previously, in order to estimate κ_{gen} in absolute error, we could stop after computing the first $\lceil 2/\varepsilon \rceil$ ergodic coefficients, but for controlling the approximation error with relative accuracy, the algorithm has to investigate on $S = \lceil c/((1 - \kappa_{\text{gen}})\varepsilon) \rceil$, $c \in \mathbb{R}_+$, which is unknown a priori. The solution is to have $\widehat{S}(X)$ depend on N_{\min} , such that the algorithm will investigate a larger space as more samples are collected. We define the estimator $\widehat{\kappa}_{\text{gen}}^{\widehat{S}}$ such that

$$\widehat{S} \doteq \lceil \sqrt{N_{\min}/d} \rceil.$$

From the triangle inequality,

$$\begin{aligned} \left| \widehat{\kappa}_{\text{gen}}^{\widehat{S}} - \kappa_{\text{gen}} \right| &\leq \left| \widehat{\kappa}_{\text{gen}}^{\widehat{S}} - \widehat{\kappa}_{\text{gen}}^{\lceil 3/(\varepsilon(1-\kappa_{\text{gen}})) \rceil} \right| \\ &\quad + \left| \widehat{\kappa}_{\text{gen}}^{\lceil 3/(\varepsilon(1-\kappa_{\text{gen}})) \rceil} - \kappa_{\text{gen}}^{\lceil 3/(\varepsilon(1-\kappa_{\text{gen}})) \rceil} \right| \\ &\quad + \left| \kappa_{\text{gen}}^{\lceil 3/(\varepsilon(1-\kappa_{\text{gen}})) \rceil} - \kappa_{\text{gen}} \right|. \end{aligned}$$

It is easy to verify that $|\kappa_{\text{gen}}^{\lceil 3/(\varepsilon(1-\kappa_{\text{gen}})) \rceil} - \kappa_{\text{gen}}| \leq \frac{(1-\kappa_{\text{gen}})\varepsilon}{3}$. To bound the second term with high-probability, we use Theorem 5 at precision $(1-\kappa_{\text{gen}})\varepsilon/3$ and confidence level $1-\delta/3$. It remains to analyze the first term:

$$\begin{aligned} & \left| \widehat{\kappa}_{\text{gen}}^{\lceil \hat{S} \rceil} - \widehat{\kappa}_{\text{gen}}^{\lceil 3/(\varepsilon(1-\kappa_{\text{gen}})) \rceil} \right| = \left| \max_{\lceil \hat{S} \rceil} \left\{ \frac{1 - \widehat{\kappa}_s}{s} \right\} - \max_{\lceil 3/((1-\kappa_{\text{gen}})\varepsilon) \rceil} \left\{ \frac{1 - \widehat{\kappa}_s}{s} \right\} \right| \\ & \leq \max \left\{ \frac{1}{s} : s \in [\hat{S} \dots \lceil 3/((1-\kappa_{\text{gen}})\varepsilon) \rceil] \cup [\lceil 3/((1-\kappa_{\text{gen}})\varepsilon) \rceil \dots \hat{S}] \right\} \\ & \leq \max \left\{ \frac{1}{\hat{S}}, \frac{(1-\kappa_{\text{gen}})\varepsilon}{3} \right\}, \end{aligned}$$

such that for $m \geq 36 \frac{d}{\pi_\star(1-\kappa_{\text{gen}})^2\varepsilon^2}$,

$$\begin{aligned} \mathbb{P}_\pi \left(\left| \widehat{\kappa}_{\text{gen}}^{\lceil \hat{S} \rceil} - \widehat{\kappa}_{\text{gen}}^{\lceil 3/(\varepsilon(1-\kappa_{\text{gen}})) \rceil} \right| > \frac{(1-\kappa_{\text{gen}})\varepsilon}{3} \right) & \leq \mathbb{P}_\pi \left(\frac{1}{\hat{S}} > \frac{(1-\kappa_{\text{gen}})\varepsilon}{3} \right) \\ & \leq \mathbb{P}_\pi \left(\sqrt{\frac{d}{N_{\min}}} > \frac{(1-\kappa_{\text{gen}})\varepsilon}{3} \right) \\ & \leq \mathbb{P}_\pi \left(|\pi_\star - \hat{\pi}_\star| > \frac{3}{4}\pi_\star \right), \end{aligned}$$

where $\hat{\pi}_\star$ is the plug-in estimator for π_\star defined in Hsu et al. (2015). Thus for $m \geq c \frac{t_{\text{mix}}}{\pi_\star} \ln \left(\frac{d}{\delta} \right)$, $c \in \mathbb{R}_+$ (Wolfer and Kontorovich, 2019, Theorem 1), this is smaller than $\delta/3$. Finally, the result is extended to non-stationary chains with Paulin (2015, Proposition 3.10). Remark: This expression for \hat{S} confirms the practical choice we proposed at (14).

4.5. Proof of Corollary 7

Combining Theorem 1 and Theorem 6, and choosing $\hat{t}_{\text{mix}} \doteq \frac{1}{1-\widehat{\kappa}_{\text{gen}}}$, for the value of m in Theorem 6, with probability at least $1-\delta$,

$$\frac{t_{\text{mix}}}{(1+\varepsilon)(1+\ln 4)} \leq \hat{t}_{\text{mix}} \leq \frac{2t_{\text{mix}}}{(1-\varepsilon)},$$

and setting $\varepsilon = 1/4$ yields the corollary. \square

5. Auxiliary facts

The following lemma and facts are proved in the appendix.

Lemma 10 *Let $X_1, \dots, X_m \sim \mathbf{M}$ stationary ergodic Markov chain. For a skipping rate s , and for $m \geq c \ln \frac{d}{\delta\varepsilon} \frac{t_{\text{mix}}}{\pi_\star}$, $c \in \mathbb{R}_+$,*

$$\mathbb{P}_\pi \left(N_{\min}^{(s)} < \frac{1}{2} \lceil (m-1)/s \rceil \pi_\star \right) \leq \frac{\delta}{2d \lceil 2/\varepsilon \rceil}.$$

Fact 5.1 For two Markov matrices \mathbf{M}_1 and \mathbf{M}_2 ,

$$|\kappa(\mathbf{M}_1) - \kappa(\mathbf{M}_2)| \leq \|\mathbf{M}_1 - \mathbf{M}_2\|_\infty.$$

Fact 5.2 Let $X_1, \dots, X_m \sim \mathbf{M}$ a stationary ergodic Markov chain with stationary distribution π and mixing time at most t_{mix} . Then the mixing time $t_{\text{mix}}^{(s)}$ of the skipped chain for $s \in [m]$,

$$X_1, X_{1+s}, X_{1+2s}, \dots, X_{1+\lfloor (m-1)/s \rfloor s} \sim (\boldsymbol{\mu}, \mathbf{M}^s),$$

is such that $t_{\text{mix}}^{(s)} \leq \lceil t_{\text{mix}}/s \rceil$.

6. Discussion and future research directions

The present work offers a new perspective on the problem of estimating the mixing properties of a Markov chain, switching the focus from spectral methods and t_{rel} to contraction methods and t_{mix} itself. This offers a first step in determining whether these two statistical problems are of equivalent complexity. The proposed algorithms are primarily of theoretical interest, as they remain computationally intensive both in space and time. Algorithmic optimization of the search over the subset S , for example by leveraging additional properties of κ_{gen} is on our research agenda.

Acknowledgments

This research was supported by the Lynn and William Frankel Center for Computer Science at Ben-Gurion University and by the Israel Science Foundation. The author thanks the anonymous referees for their valuable comments.

References

- Pierre Brémaud. *Markov Chains: Gibbs Fields, Monte Carlo Simulation, and Queues*. Springer-Verlag, 1999.
- Russ Bubley and Martin Dyer. Path coupling: A technique for proving rapid mixing in markov chains. In *Proceedings 38th Annual Symposium on Foundations of Computer Science*, pages 223–231. IEEE, 1997.
- Kai-Min Chung, Henry Lam, Zhenming Liu, and Michael Mitzenmacher. Chernoff-hoeffding bounds for markov chains: Generalized and simplified. *arXiv preprint arXiv:1201.0559*, 2012.
- Richard Combes and Mikael Touati. Computationally efficient estimation of the spectral gap of a markov chain. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 3(1):7, 2019.
- RL Dobrushin. Central limit theorem for non-stationary markov chains, i, ii. *theory prob. appl.* 1, 65-80, 329-383. *English translation*, 1956.
- Martin Dyer, Leslie Ann Goldberg, Catherine Greenhill, Mark Jerrum, and Michael Mitzenmacher. An extension of path coupling and its application to the glauber dynamics for graph colorings. *SIAM Journal on Computing*, 30(6):1962–1975, 2001.

- Daniel Hsu, Aryeh Kontorovich, David A. Levin, Yuval Peres, Csaba Szepesvri, and Geoffrey Wolfer. Mixing time estimation in reversible markov chains from a single sample path. *Ann. Appl. Probab.*, 29(4):2439–2480, 08 2019. doi: 10.1214/18-AAP1457. URL <https://doi.org/10.1214/18-AAP1457>.
- Daniel J Hsu, Aryeh Kontorovich, and Csaba Szepesvári. Mixing time estimation in reversible markov chains from a single sample path. In *Advances in neural information processing systems*, pages 1459–1467, 2015.
- Daniel Jerison. General mixing time bounds for finite markov chains via the absolute spectral gap. *arXiv preprint arXiv:1310.8021*, 2013.
- Jiantao Jiao, Yanjun Han, and Tsachy Weissman. Minimax estimation of the $l_{\{1\}}$ distance. *IEEE Transactions on Information Theory*, 64(10):6672–6706, 2018.
- Sudeep Kamath and Sergio Verdú. Estimation of entropy rate and rényi entropy rate for markov chains. In *2016 IEEE International Symposium on Information Theory (ISIT)*, pages 685–689. IEEE, 2016.
- David A Levin and Yuval Peres. Estimating the spectral gap of a reversible markov chain from a short trajectory. *arXiv preprint arXiv:1612.05330*, 2016.
- David Asher Levin, Yuval Peres, and Elizabeth Lee Wilmer. *Markov chains and mixing times, second edition*. American Mathematical Soc., 2009.
- Malwina Luczak. Concentration of measure and mixing for markov chains. *Discrete Mathematics & Theoretical Computer Science*, 2008.
- Yann Ollivier. Ricci curvature of markov chains on metric spaces. *Journal of Functional Analysis*, 256:810–864, 2009.
- Daniel Paulin. Concentration inequalities for Markov chains by Marton couplings and spectral methods. *Electronic Journal of Probability*, 20, 2015.
- Daniel Paulin. Mixing and concentration by ricci curvature. *Journal of Functional Analysis*, 270(5):1623–1662, 2016.
- Qian Qin, James P Hobert, Kshitij Khare, et al. Estimating the spectral gap of a trace-class markov operator. *Electronic Journal of Statistics*, 13(1):1790–1822, 2019.
- Geoffrey Wolfer and Aryeh Kontorovich. Estimating the mixing time of ergodic markov chains. In *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 3120–3159. PMLR, 2019. URL <http://proceedings.mlr.press/v99/wolfer19a.html>.

Proof of auxiliary lemmas and facts

Lemma 10 Let $X_1, \dots, X_m \sim M$ stationary ergodic Markov chain. For a skipping rate s , and for $m \geq c \ln \frac{d}{\delta \varepsilon} \frac{t_{\text{mix}}}{\pi_\star}$, $c \in \mathbb{R}_+$,

$$\mathbb{P}_\pi \left(N_{\min}^{(s)} < \frac{1}{2} \lceil (m-1)/s \rceil \pi_\star \right) \leq \frac{\delta}{2d \lceil 2/\varepsilon \rceil}.$$

Proof From the Chernoff-Hoeffding lower tail concentration inequality at [Chung et al. \(2012, Theorem 3.1\)](#), for $X_1, \dots, X_m \sim M$ ergodic over d states, stationary, with mixing time $t_{\text{mix}}(\xi)$ for $\xi \leq 1/8$, and $\eta \in (0, 1)$,

$$\mathbb{P}_\pi (N_i \leq (1-\eta)(m-1)\pi(i)) \leq c \exp \left(-\frac{\eta^2 m \pi(i)}{72 t_{\text{mix}}(\xi)} \right), c \in \mathbb{R}_+$$

where we already used the definition of N_i and that by stationarity, $\mathbb{E}_\pi [\mathbf{1}\{X_t = i\}] = \pi(i)$. The astute reader will notice that our definition of t_{mix} is for $\xi = 1/4$, such that this theorem is not applicable verbatim, however [Chung et al. \(2012\)](#) mentions (at p.3) that it generally holds with $\frac{1-\sqrt{2\xi}}{36 t_{\text{mix}}(\xi)}$ instead of $\frac{1}{72 t_{\text{mix}}(\xi)}$. In our case, we therefore adapt the constant to $c' = \frac{1-\sqrt{1/2}}{36}$, and for $\eta = 1/2$,

$$\mathbb{P}_\pi \left(N_i \leq \frac{(m-1)\pi(i)}{2} \right) \leq c \exp \left(-\frac{c'(m-1)\pi(i)}{t_{\text{mix}}} \right). \quad (15)$$

We proceed and apply the above to the different multi-step chains for $s \in [S]$. The mixing time $t_{\text{mix}}^{(s)}$ of each such chain is at most about $\frac{t_{\text{mix}}}{s}$, which is formalized in [Fact 5.2](#).

$$\begin{aligned} \mathbb{P}_\pi \left(N_i^{(s)} \leq \frac{1}{2} \lceil (m-1)/s \rceil \pi(i) \right) &\stackrel{(i)}{\leq} c \exp \left(-\frac{c' \lceil (m-1)/s \rceil \pi(i)}{t_{\text{mix}}^{(s)}} \right) \\ &\stackrel{(ii)}{\leq} c \exp \left(-\frac{c'(m-s-1)\pi(i)}{t_{\text{mix}}} \right) \end{aligned}$$

where (i) is (15) applied to the skipped chain, as it is still the case that $\mathbb{E}_\pi [\mathbf{1}\{X_{1+s(t-1)} = i\}] = \pi(i)$, and (ii) is [Fact 5.2](#). As a consequence for $m \geq c'' \left(\frac{1}{\varepsilon} + \frac{t_{\text{mix}}}{\pi(i)} \ln \frac{2d \lceil 2/\varepsilon \rceil}{\delta} \right)$, this error probability is smaller than $\frac{\delta}{2d \lceil 2/\varepsilon \rceil}$. Taking a maximum over $i \in \Omega$, yields the lemma. \blacksquare

Fact 5.1 For two Markov matrices M_1 and M_2 ,

$$|\kappa(M_1) - \kappa(M_2)| \leq \|M_1 - M_2\|_\infty.$$

Proof This is a direct consequence of the sub-additivity of the sup norm.

$$\begin{aligned} 2 |\kappa(M_1) - \kappa(M_2)| &\leq \max_{(i,j) \in \Omega^2} \left| \|M_1(i, \cdot) - M_1(j, \cdot)\|_1 - \|M_2(i, \cdot) - M_2(j, \cdot)\|_1 \right| \\ &\leq \max_{(i,j) \in \Omega^2} \|M_1(i, \cdot) - M_1(j, \cdot) - M_2(i, \cdot) + M_2(j, \cdot)\|_1 \\ &\leq \max_{(i,j) \in \Omega^2} (\|M_1(i, \cdot) - M_2(i, \cdot)\|_1 + \|M_1(j, \cdot) - M_2(j, \cdot)\|_1) \\ &= 2 \max_{i \in \Omega} \|M_1(i, \cdot) - M_2(i, \cdot)\|_1 \\ &= 2 \|M_1 - M_2\|_\infty. \end{aligned}$$

■

Fact 5.2 Let $X_1, \dots, X_m \sim M$ a stationary ergodic Markov chain with stationary distribution π and mixing time at most t_{mix} . Then the mixing time $t_{\text{mix}}^{(s)}$ of the skipped chain for $s \in [m]$,

$$X_1, X_{1+s}, X_{1+2s}, \dots, X_{1+\lfloor (m-1)/s \rfloor s} \sim (\mu, M^s),$$

is such that $t_{\text{mix}}^{(s)} \leq \lceil t_{\text{mix}}/s \rceil$.

Proof Let t such that $t > \lceil t_{\text{mix}}/s \rceil$, then

$$\|\mu(M^s)^t - \pi\|_{\text{TV}} \leq \|\mu(M^s)^{\lceil t_{\text{mix}}/s \rceil} - \pi\|_{\text{TV}} = \|\mu M^{t_{\text{mix}}} - \pi\|_{\text{TV}} \leq \frac{1}{4},$$

where the first inequality holds as advancing the chain can only move it closer to stationarity (Levin et al., 2009, Exercise 4.2). ■

Algorithm We describe in Algorithm 1 the adaptive version of the procedure that outputs an estimator for the mixing time \hat{t}_{mix} , with the guarantees of Corollary 7, modulo a smoothing parameter λ .

The time complexity of the algorithm of Hsu et al. (2019) for estimating the absolute spectral gap of a reversible chain is of the order of $\mathcal{O}(m + d^3)$. The extension of Wolfer and Kontorovich (2019) that involves the first $S \in \mathbb{N}$ multiplicative reversiblizations of the chain, without considering any form of parallelization, has a computational complexity upper bounded by $\mathcal{O}(S(m + d^3))$. Algorithm 1 has an equivalent time complexity, as computing the Dobrushin contraction coefficient requires $\mathcal{O}(d^3)$.

Interestingly, the complexity of constructing the confidence interval compares favorably with the previous methods. In Hsu et al. (2019) the necessity of computing the *pseudo-inverse* of the empirical transition matrix leads to a complexity of $\mathcal{O}(m + d^3)$. In Wolfer and Kontorovich (2019) computing an interval requires $\tilde{\mathcal{O}}(m + d^2)$ for a reversible chain, and $\mathcal{O}(S(m + d^3))$ in the non-reversible case. In our algorithm, computing the interval can be done in $\mathcal{O}(S(m + d))$.

```

Function MixingTime ( $d, \lambda, (X_1, \dots, X_m)$ ):
    | return  $1 / (1 - \text{GeneralizedContractionCoeffAdaptive}(d, \lambda, (X_1, \dots, X_m)))$ 
Function GeneralizedContractionCoeffAdaptive ( $d, \lambda, (X_1, \dots, X_m)$ ):
    |  $\mathbf{N} \leftarrow [0]_d$       $N_{\min} \leftarrow m$ 
    | for  $t \leftarrow 1$  to  $m - 1$  do
    | |  $\mathbf{N}[X_t] \leftarrow \mathbf{N}[X_t] + 1$ 
    | |
    | | end
    | | for  $i \leftarrow 1$  to  $d$  do
    | | | if  $\mathbf{N}[i] < N_{\min}$  then
    | | | |  $N_{\min} \leftarrow \mathbf{N}[i]$ 
    | | | | end
    | | | end
    | | end
    | | return  $\text{GeneralizedContractionCoeff}(d, \lambda, (X_1, \dots, X_m), \lceil \sqrt{N_{\min}/d} \rceil)$ 
Function GeneralizedContractionCoeff ( $d, \lambda, (X_1, \dots, X_m), S$ ):
    |  $r_{\max} \leftarrow 0$ 
    | for  $s \leftarrow 1$  to  $S$  do
    | |  $\kappa \leftarrow \text{ContractionCoeff}(d, \lambda, (X_1, X_{1+s}, X_{1+2s}, \dots, X_{1+\lfloor (m-1)/s \rfloor s}))$ 
    | | if  $(1 - \kappa)/s > r_{\max}$  then
    | | |  $r_{\max} \leftarrow (1 - \kappa)/s$ 
    | | | end
    | | end
    | | end
    | | return  $1 - r_{\max}$ 
Function ContractionCoeff ( $d, \lambda, (X_1, \dots, X_n)$ ):
    |  $\mathbf{N} \leftarrow [d\lambda]_d$ 
    |  $\mathbf{T} \leftarrow [\lambda]_{d \times d}$ 
    | for  $t \leftarrow 1$  to  $n - 1$  do
    | |  $\mathbf{N}[X_t] \leftarrow \mathbf{N}[X_t] + 1$ 
    | |  $\mathbf{T}[X_t, X_{t+1}] \leftarrow \mathbf{T}[X_t, X_{t+1}] + 1$ 
    | |
    | | end
    | |  $a_{\max} \leftarrow 0$ 
    | | for  $i \leftarrow 1$  to  $d$  do
    | | | for  $j \leftarrow 1$  to  $d$  do
    | | | |  $a = 0$ 
    | | | | for  $k \leftarrow 1$  to  $d$  do
    | | | | |  $a \leftarrow a + |\mathbf{T}[i, k]/\mathbf{N}[i] - \mathbf{T}[j, k]/\mathbf{N}[j]|$ 
    | | | | | end
    | | | | end
    | | | | if  $a > a_{\max}$  then
    | | | | |  $a_{\max} \leftarrow a$ 
    | | | | | end
    | | | | end
    | | | end
    | | | end
    | | | return  $a_{\max}/2$ 

```

Algorithm 1: The estimation procedure outputting \hat{t}_{mix} .